

WHAT IS CLAIMED IS:

1. A fault-tolerant computer network file system, comprising:  
a first file server operably connected to a network fabric;  
a second file server operably connected to the network fabric;  
a first disk array operably coupled to said first file server and to said  
second file server;

5 a second disk array operably coupled to said first file server and to said  
second file server;

10 first file system information loaded on said first file server, said first file  
system information comprising a first intent log of proposed changes to first  
metadata;

15 second file system information loaded on said second file server, said  
second file system information comprising a second intent log of proposed  
changes to second metadata, said first file server having a copy of said second  
metadata, said second file server maintaining a copy of said first metadata,  
thereby allowing said first file server to access files on said second disk array in  
the event of a failure of said second file server.

20 2. The fault-tolerant computer network file system of Claim 1, wherein  
said first metadata further comprises first directory information that describes a  
directory structure of a portion of the network file system whose directories are stored  
on said first file server and second directory information that describes a directory  
structure of a portion of the network file system whose directories are stored on said  
second file server, said first directory information comprising location information for  
said second directory information, said location information comprising a server id that  
identifies said second file server.

25 3. The fault-tolerant computer network file system of Claim 1, wherein said  
first metadata comprises a root directory.

30 4. The fault-tolerant computer network file system of Claim 1, wherein said  
first metadata further comprises directory information that describes a directory  
structure of a portion of the network file system whose directories are stored on said  
first file server, said directory information comprising location information for a first

- PCT/US2013/050000
- file, said location information comprising a server id that identifies at least said first file server or said second file server.
5. The fault-tolerant computer network file system of Claim 4, wherein said location information further comprises a disk id that identifies a disk drive.
- 5 6. The fault-tolerant computer network file system of Claim 4, wherein said directory information comprises at least one Gnid-string.
7. The fault-tolerant computer network file system of Claim 6, wherein a one-to-one correspondence exists between said at least one Gnid-string and a directory stored on said first file server.
- 10 8. The fault-tolerant computer network file system of Claim 6, wherein said Gnid-string comprises a collection of gnids.
9. The fault-tolerant computer network file system of Claim 8, wherein each of said gnids comprises information for locating a specified gnode.
- 15 10. The fault-tolerant computer network file system of Claim 9, wherein said information for locating a specified gnode comprises a pointer to said specified gnode.
11. The fault-tolerant computer network file system of Claim 9, wherein said specified gnode comprises file attributes for a file corresponding to said gnode.
- 20 12. The fault-tolerant computer network file system of Claim 11, wherein said file attributes include at least one of a file id, a file access time, a file creation time, and a file modification time.
13. The fault-tolerant computer network file system of Claim 9, wherein said specified gnode comprises information for locating a first gee of a plurality of gees corresponding to said gnode.
- 25 14. The fault-tolerant computer network file system of Claim 13, wherein said plurality of gees comprises gnode gees and data gees.
15. The fault-tolerant computer network file system of Claim 14, wherein each of said gnode gees comprises information to specify an extent.
- 30 16. The fault-tolerant computer network file system of Claim 14, wherein each of said data gees comprises information to specify a first logical disk block and information to specify a disk that contains said first logical block.

17. The fault-tolerant computer network file system of Claim 14, wherein  
said plurality of gees further comprises parity gees.
18. The fault-tolerant computer network file system of Claim 17, wherein  
each of said parity gees comprises information regarding location of parity data for one  
or more preceding data gees in said plurality of gees.  
5
19. The fault-tolerant computer network file system of Claim 14, wherein a  
parity group comprises a first set of one or more data gees and an associated parity gee.
20. The fault-tolerant computer network file system of Claim 19, wherein  
each data gee identifies a block of data and said parity gee identifies a parity block.  
10
21. The fault-tolerant computer network file system of Claim 20, wherein  
each block of data and parity in said parity group is stored on a separate disk drive such  
that no single disk drive contains data from two blocks of said parity group.
22. The fault-tolerant computer network file system of Claim 19, wherein a  
size of a first parity group is independent of a size of a second parity group.  
15
23. The fault-tolerant computer network file system of Claim 4, wherein said  
first directory information is mirrored on said second file server.
24. The fault-tolerant computer network file system of Claim 1, wherein said  
network fabric comprises a Fibre channel network.  
20
25. The fault-tolerant computer network file system of Claim 1, wherein said  
network fabric comprises an ethernet network.
26. The fault-tolerant computer network file system of Claim 1, wherein said  
network fabric comprises an asynchronous transfer mode network.  
25
27. The fault-tolerant computer network file system of Claim 1, wherein said  
network fabric comprises a first Fibre channel network and wherein said first file server  
communicates with said first disk array and said second disk array using a second Fibre  
channel network.
28. The fault-tolerant computer network file system of Claim 1, wherein files  
stored on said first disk array and files stored on said second disk array are located in a  
hierarchical directory structure having a common root directory, said first file system  
information comprising directory information that describes directories stored on said  
first file server, said second file system information comprising directory information  
30

that describes directories stored on said second file server, said first directory information comprising location information for finding said second directory information, said location information comprising a server id that identifies said second file server.

5        29. The fault-tolerant computer network file system of Claim 1, wherein files stored by said first disk array and files stored by said disk array are located in a hierarchical directory structure having a common root directory, said first file system information comprising metadata for locating files stored on said second disk array.

10      30. The fault-tolerant computer network file system of Claim 1, wherein said first file system information comprises first metadata that describes directories stored on said first disk array, said file system information comprising second metadata that describes directories stored on said disk array, said first metadata comprising location information for locating said second metadata, said location information comprising a server id.

15      31. The fault-tolerant computer network file system of Claim 30, wherein said first metadata comprises a root directory.

32. The fault-tolerant computer network file system of Claim 30, wherein at least said first metadata comprises file attributes for one or more files stored by said first file server.

20      33. The fault-tolerant computer network file system of Claim 30, wherein at least said first metadata comprises information to specify a selected logical disk block of a selected file and information to specify a disk that contains said selected logical block.

25      34. The fault-tolerant computer network file system of Claim 30, wherein said metadata identifies data blocks and parity blocks corresponding to said data blocks.

30      35. The fault-tolerant computer network file system of Claim 30, wherein said metadata identifies parity groups, said parity groups comprising a plurality of information blocks, said information blocks comprising one or more data blocks, said information blocks further comprising a parity block, each of said information blocks stored on a different disk drive.

36. The fault-tolerant computer network file system of Claim 1, wherein a size of a first parity group is independent of a size of a second parity group.
- 5 37. The fault-tolerant computer network file system of Claim 1, wherein said copy of said first intent log is configured to provide sufficient information to allow said second file server to modify and access files on said first disk array without interruption in the event of a failure of said first file server.
- 10 38. The fault-tolerant computer network file system of Claim 1, wherein said copy of said first intent log is configured to provide sufficient information to allow said first file server to be hot-swapped without loss of information stored on said first disk array.
- 15 39. The computer network file system of Claim 1, wherein said first file server communicates with said first disk array using a Fibre channel network.
40. The computer network file system of Claim 1, wherein said first file server communicates with said first disk array using InfiniBand.
- 15 41. The computer network file system of Claim 1, wherein said first file server communicates with said first disk array using SCSI.
42. The computer network file system of Claim 1, wherein said first file server receives a copy of entries in said second intent log and said second file server receives a copy of entries in said first intent log.
- 20 43. The computer network file system of Claim 1, wherein said first file server receives a copy of entries in said second intent log and said second file server receives a copy of entries in said first intent log.
- 25 44. The computer network file system of Claim 1, wherein said first file server is configured to detect when said second file server goes offline and begin servicing file requests for files owned primarily by said second file server.
45. The computer network file system of Claim 44, wherein said first file server is configured to detect when said second file server goes online after being offline, said first file server configured to allow said second file server to service file requests for files owned primarily by said second server.
- 30 46. A method for hot-swapping file servers in a computer network, comprising:

loading first file system metadata on a first file server operably connected to a network fabric, said first file system operably connected to a first disk drive and a second disk drive;

5 loading second file system metadata on a second file server connected to said network fabric, said second file system operably connected to said first disk drive and to said second disk drive;

copying a first intent log from said first file server to a backup intent log on said second file server, said first intent log providing information regarding future changes to information stored on said first disk drive; and

10 using said backup intent log to allow said second file server to make changes to said information stored on said first disk drive.

47. The method of Claim 46, further comprising storing first file system directory information on said first file server, said first file system directory information describing a directory structure of a portion of the network file system whose directories are stored on said first file server, said first file system directory information comprising location information for a first file, said location information comprising a server id that identifies said second file server.

20 48. The method of Claim 46, further comprising storing first file system directory information on said first file server, said first file system directory information describing a directory structure of a portion of the network file system whose directories are stored on said first file server, said first intent log comprising changes to be made to said first file system directory information.

49. The method of Claim 48, wherein said first file system directory information comprises a root directory.

25 50. The method of Claim 48, wherein said file system directory information comprises at least one Gnid-string.

51. The method of Claim 50, wherein a correspondence exists between said at least one Gnid-string and a directory stored on said first file server.

30 52. The method of Claim 50, wherein said Gnid-string comprises a collection of gnids.

- 2016 RELEASE UNDER E.O. 14176
53. The method of Claim 52, wherein each of said gnids comprises information for locating a specified gnode.
54. The method of Claim 50, wherein said information for locating a specified gnode comprises a pointer to said specified gnode.
- 5 55. The method of Claim 54, wherein said specified gnode comprises file attributes for a file corresponding to said gnode.
56. The method of Claim 55, wherein said file attributes include at least one of a file id, a file access time, a file creation time, and a file modification time.
- 10 57. The method of Claim 54, wherein said specified gnode comprises information for locating a first gee of a plurality of gees corresponding to said gnode.
58. The method of Claim 57, wherein said plurality of gees comprises gnode gees and data gees.
- 15 59. The method of Claim 58, wherein each of said gnode gees comprises information to specify a logical block extent.
60. The method of Claim 58, wherein each of said data gees comprises information to specify a first logical disk block and information to specify a disk that contains said first logical block.
- 15 61. The method of Claim 58, further comprising defining a plurality of parity gees.
- 20 62. The method of Claim 61, wherein each of said parity gees comprises information regarding location of parity data for one or more preceding data gees in said plurality of gees.
- 25 63. The method of Claim 46, further comprising defining at least one parity group having a first parity group size, said at least one parity group comprising a parity block and one or more data blocks.
64. The method of Claim 63, further comprising storing each of said data blocks and said parity block on different disk drives.
65. The method of Claim 63, further comprising defining at least one parity group having a second parity group size.

66. The method of Claim 46, further comprising storing file system directory information that describes at least a portion of a hierarchical directory structure, said hierarchical directory structure spanning said first file server and said second file server.

5 67. The method of Claim 46, wherein said first file server continues operation of said second file server in the event of a failure of said second file server.

68. A computer network file system, comprising:

a first file server operably connected to a network fabric and to a first disk array and to a second disk array;

10 a second file server operably connected to said network fabric and to said first disk array and to said second disk array;

means for locating files stored by said first file server and files stored by said second file server by traversing a directory structure that spans at least said first file server and said second file server; and

15 means for allowing said first file server to complete file system changes intended by said second file server but not completed due to said second file server going offline.

69. The computer network file system of Claim 68, wherein said directory structure comprises location information for a first file, said location information comprising a server id that identifies at least said first file server or said second file server.

20 70. The computer network file system of Claim 68, wherein said directory structure comprises server ids of servers that contain sub-directories.

71. The computer network file system of Claim 68, further comprising means for sharing workload between said first file server and said second file server.

25 72. The computer network file system of Claim 68, further comprising means for detecting that said second file server has come online and handing at least a portion of file system operations dealing with said second disk array over to said second file server.